

School of Computing and Information Systems  
The University of Melbourne  
COMP90049 Knowledge Technologies (Semester 1, 2019)  
Workshop exercises: Week 6

1. In the context of an **information retrieval** engine, what does it mean for a document to be **returned** for a query? What does it mean for a document to be **relevant** for a query?
2. Identify some differences between **Boolean** querying and **ranked** querying, in an Information Retrieval context.
  - (a) (Extension) Obviously, search engines like Google use ranked querying by default ... or do they? What evidence do we have that there is a “hybrid” Boolean component in Google’s typical behaviour?
3. Identify the two (sometimes three) components of a **TF-IDF model**. Indicate the rationale behind them as in, why would they contribute to a “better” result set?
4. Many TF-IDF models are possible; consider the following one:

$$w_{d,t} = \begin{cases} 1 + \log_2 f_{d,t} & \text{if } f_{d,t} > 0 \\ 0 & \text{otherwise} \end{cases}$$
$$w_{q,t} = \begin{cases} \log(1 + \frac{N}{f_t}) & \text{if } f_{q,t} > 0 \\ 0 & \text{otherwise} \end{cases}$$

- (a) Construct suitable vectors for the five documents in the collection below, and then use the **cosine measure** to determine the document ranking for the (conjunctive) query **apple lemon**:

<i>DocID</i>	<b>apple</b>	<b>ibm</b>	<b>lemon</b>	<b>sun</b>
Doc <sub>1</sub>	4	0	0	1
Doc <sub>2</sub>	5	0	5	0
Doc <sub>3</sub>	2	5	0	0
Doc <sub>4</sub>	1	2	1	7
Doc <sub>5</sub>	1	1	3	0

- (b) If Documents 4 and 5 were the only **truly** relevant documents in the collection, calculate  $P@1$ ,  $P@3$ , and  $P@5$  for the above system.
5. What does **Recall** correspond to, in an Information Retrieval context? Why is Recall not usually considered to be a useful evaluation metric for an IR Engine? How does this relate to  $P@k$ ?